

STARFISH

take the “meh” out of metadata

harness the “chi” in archiving

take the “rage” out of storage

put the “tada” in metadata

What Does Starfish Do?

- Starfish associates metadata with files and directories in POSIX-style file systems and in S3-style object stores.
- Metadata are used for:
 - Content classification
 - Reporting and analytics
 - User portal
 - Shaping policies and directing batch operations
- There are two forms of metadata
 - Simple tags (yes, directory tags inherit)
 - Key-value pairs (using JSON)
- Starfish sits outside of the data path so it must asynchronously maintain the inventory of the file system.
 - It works with all storage systems
 - It does not introduce any latency or points of failure
 - It is very fast and efficient
 - Starfish is suitable for extreme environments with billions of files and 100s of petabytes.

Two Main Use Cases Among Librarians and Curators

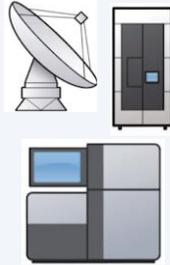
1. A middleware for storage housekeeping in a data curation facility.
 - Data protection
 - High-level content classification
 - Fixity
 - Reporting, cost-accounting, charge-back
 2. A unique solution for associating metadata and automating workflows with live, mutable file storage systems.
 - Allows the curation workflow to begin with data creation
 - Facilitates the hand-off between content creation and curation
 - Facilitates the re-use of archival data sets
- Starfish is also great for managing files that are awaiting formal curation workflows.

Life Without Starfish

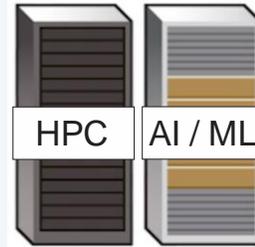
Applications that Reference Files

Content Repositories
Digital Asset Management
Data Commons
Lab Notebooks

Data Acquisition Workflows



Data Processing Workflows



Tribal Knowledge (or lack thereof)



HPC Storage



Enterprise NAS



Tape Archive

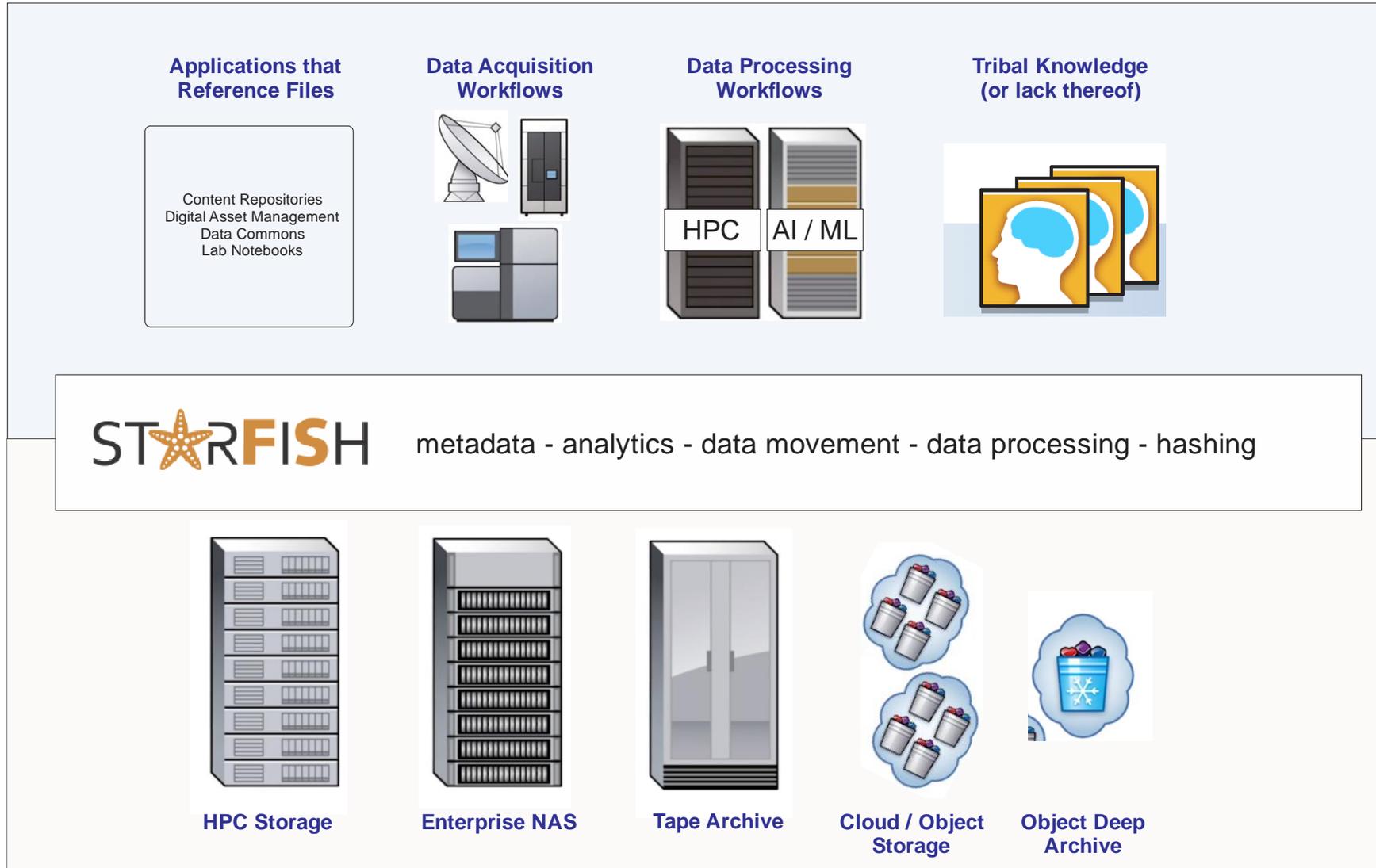


Cloud / Object Storage



Object Deep Archive

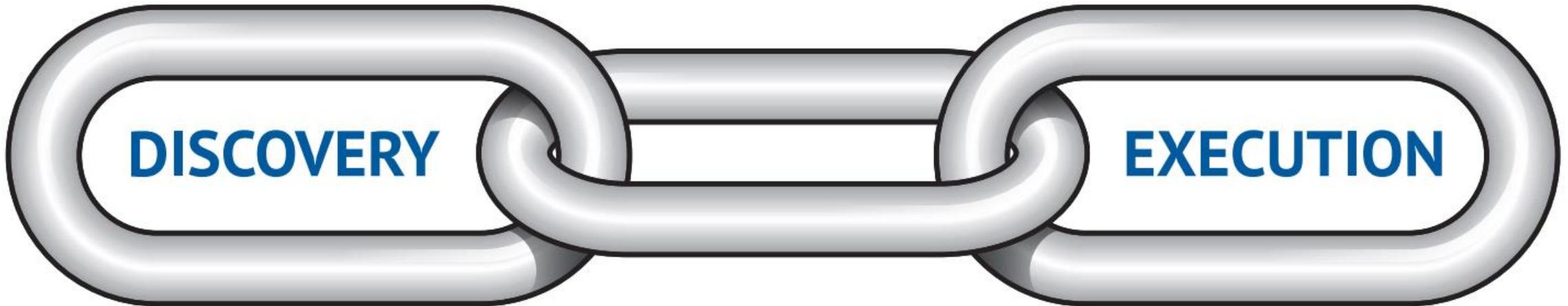
Starfish is as Storage Middleware



A Simpler Way to Explain the Technology

*If your files could talk,
what could they tell you
about themselves?*

*If your files could listen
and obey, what would you
tell them to do?*



Discovery + Execution: A Simple But Powerful Paradigm

*If your files could talk,
what could they tell you
about themselves?*

DISCOVERY

a data catalog for unstructured data
massively scalable - billions of files and objects
extensible metadata - tags and key-value pairs
simple and turnkey but suitable for custom integration

*If your files could listen
and obey, what would you
tell them to do?*

EXECUTION

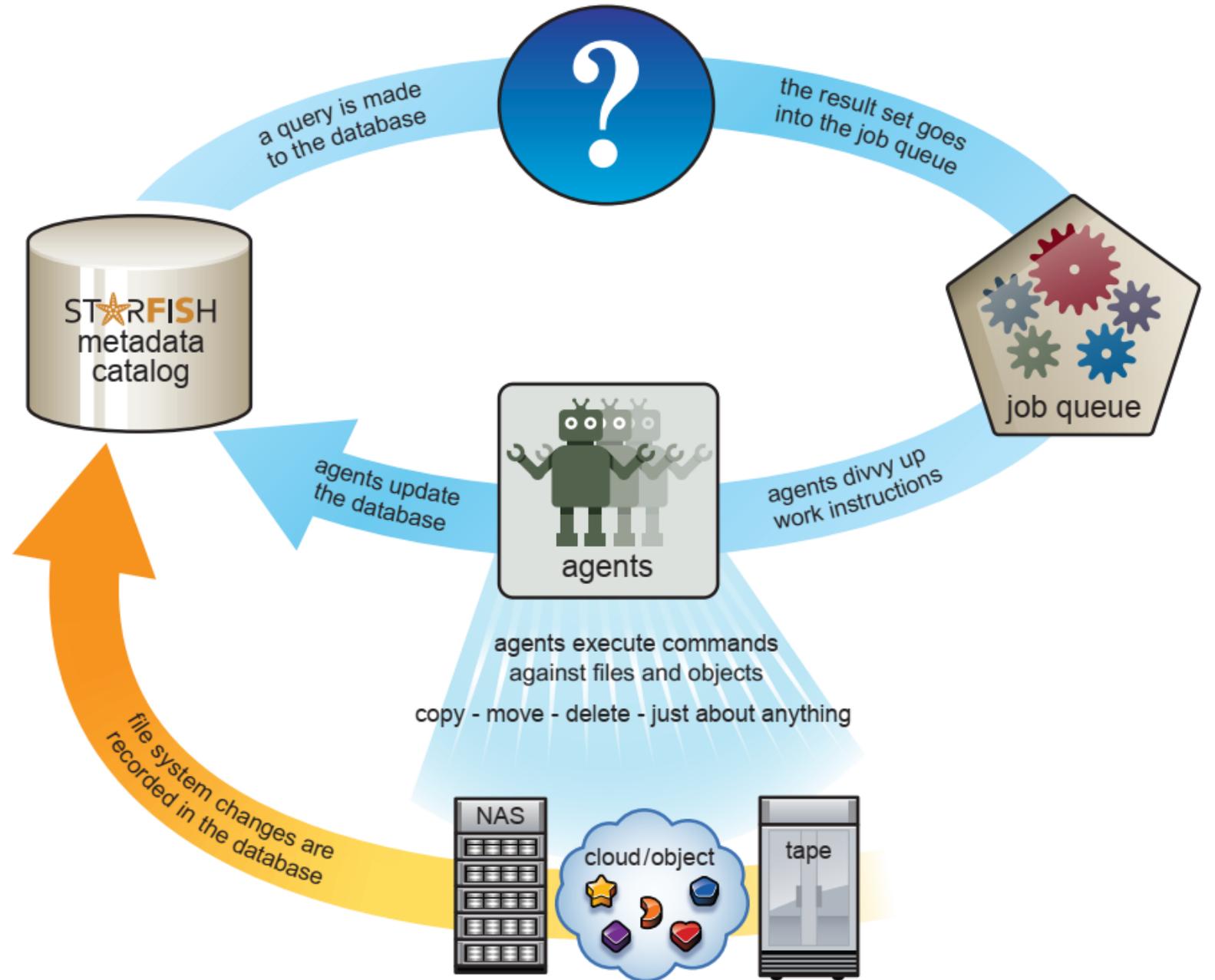
a scale-out data mover & batch processor
LINUX or Windows – File and Object
migrate - replicate - workflow - archive - backup - purge
easy to customize - runs your code or ours

Lather, Rinse, Repeat

There is a feedback loop between the catalog and the jobs engine.

- Query database
- Do some stuff
- Update database

- Query the database
- Do more stuff
- Update database



Key-value metadata is simple JSON

STARFISH 6.0.7537

Analytics Browser Tags Zones Jobs Scans

HELP HINTS SETTINGS

VOLUMES

Logical size Physical size

weka U: 12.98 GiB; F: 214.39 GiB

windows U: 1.22 MiB; F: 227.36 GiB

Location: chicago

netapp U: 1.42 GiB; F: 225.94 GiB

Location: cloud

amazon U: 4 KiB; F: 227.37 GiB

azure U: 4 KiB; F: 89.56 GiB

s3-bigbucket U: 10.51 GiB; F: 216.86 GiB

sf-bigbucket U: 735.99 MiB; F: 1 PiB

Edit Condensed view 17

ZONES

TAG SEARCH

SEARCH

DATA AS OF

prod

Name	Rec...	Rec...	Log...	Use...	Cou...	Cost
> material...	2019-0...	2020-0...	53.33 ...	fred	288	0.01
> material...	2018-0...	2020-0...	20 KiB	root	4	0.00
> material...	2018-0...	2020-0...	373.3 ...	betty	2,128	0.04
> medical ...	2021-0...	2021-0...	28.17 ...	john	187,449	3,097.53
> mthead ...	2019-0...	2020-0...	63.28 ...	root	383	0.01
> Nthead ...	2019-0...	2019-0...	56.64 ...	root	232	0.01
> oldimag...	2021-0...	2020-1...	100 TiB	richard	477,165	10,995.25
> Othead ...	2019-1...	2020-0...	1.2 GiB	root	6,869	0.13
> Pictures ...	2021-0...	2021-0...	1.89 GiB	demo	2,869	0.20
> 23 ...	-	-	45.14 ...		23	0.00
000	-	-	1.62 M...	demo	1	0.00
000	-	-	564 KiB	demo	1	0.00
000	-	-	547.36...	demo	1	0.00
Summary	-	-	144.84...		1,509,...	15,925.43

Job - LookforPii

Has Pii	False
Filename	/vois/production/Pictures/0002.DCM
Scan End	2020-10-09T20:39:31.925798+00:00
Scan Start	2020-10-09T20:39:31.925794+00:00
Supported	No
Time Executed	2020-10-09 16:39

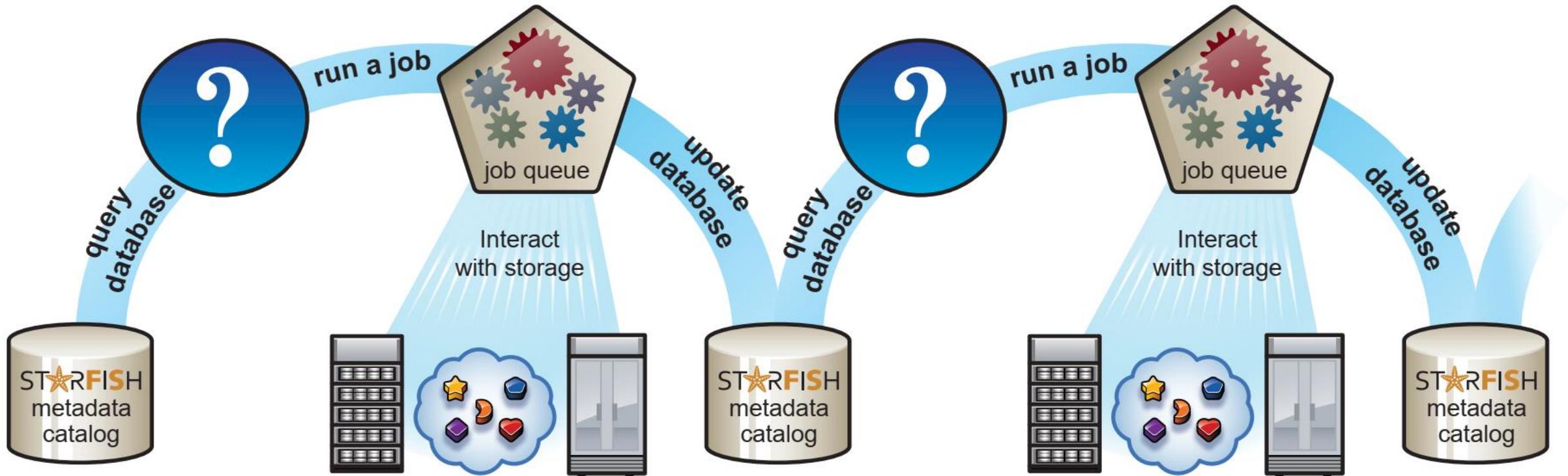
Job - hash

md5	bdc98d68248228bf79370cb6807803b
sha1	08efb06808fd4ee5d3a61407ed5e914f
Time Executed	2020-10-09 16:39

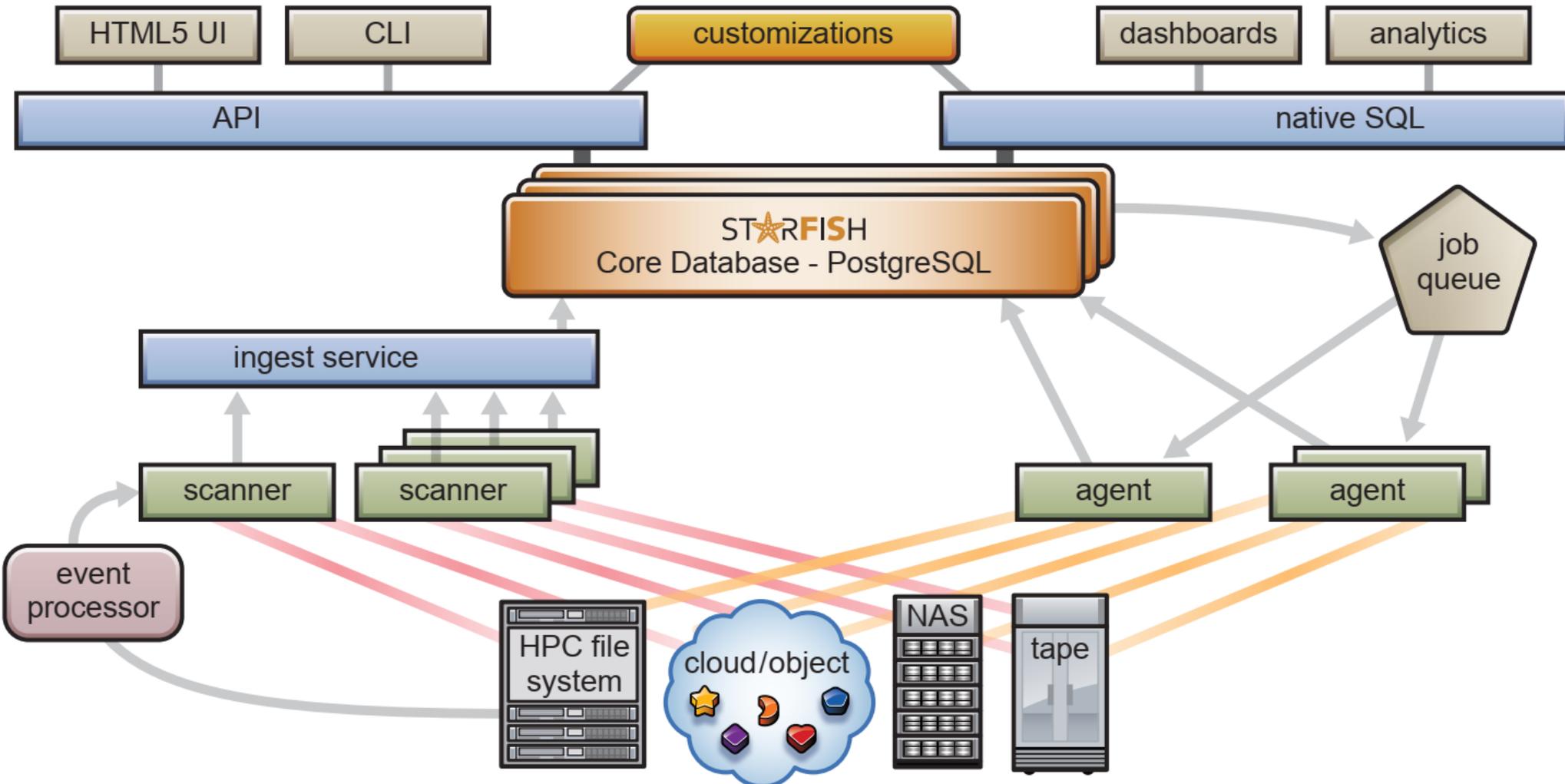
Job - meta-x

Accession Number	
Axis Units	[DPPS, 'NONE']
Bits Allocated	8
Bits Stored	8
Columns	512
Coordinate Start Value	0
Coordinate Step Value	40
Curve Data Descriptor	[0, 1]
Curve Dimensions	2
Curve Range	

Knowledge Begets Knowledge



Starfish Topology: Designed for Scale and Customization



The Grand Vision

Governance & Stewardship

Archivists - Curators - Librarians - Knowledge Managers
Data Governors - Auditors - Data Wranglers

Retention / Disposition
Anonymization / Privacy
Compliance (GDPR - Export Controls. . . .)
Data integrity assurance
Content classification
Data provenance
Data cataloging



Creation and Consumption

Researchers - Artists - Animators - Engineers - Quants
Core Facilities - Data Scientists

Metadata tagging
Search
Workflow automation
Data provenance
Collaboration
Data reusability
Data transmission

Sys Admins (Storage, Backup, HPC) - IT Directors

Data Movement: Archiving/Tiering, Migration, Replication, Cloud-bursting, Backup & Restore

Reporting: Capacity Planning, Aging, Cost Accounting, Charge-back / Show-back

Misc Operations: Permissions Management, GDPR Compliance, Job Scheduling

Starfish is Open, Low Risk and Built for Science and R&D

- Non-proprietary
 - No proprietary file formats
 - Industry standard compression (GZIP and LZ)
 - TAR when group files together
 - Storage device agnostic
 - Metadata is stored in SQL (Postgres)
- Starfish is out of the data path
 - It does not introduce latency
 - It does not introduce points of failure
- Starfish was built for Research Data Management
 - Our long term roadmap addresses the coming challenges of research computing